

Code Obfuscation

David Pichardie

Master recherche Université Rennes 1

December 2011

Module SDL

Lecturer: Thomas Jensen

Code Obfuscation¹

1. Most of these slides are taken from C. Collberg, C. Thomborson, D. Low, *A Taxonomy of Obfuscating Transformations*, <http://www.cs.arizona.edu/~collberg/Research/Publications/CollbergThomborsonLow97a/index.html>

Motivations

Software developers (Alice) want to protect their code from reverse engineering.

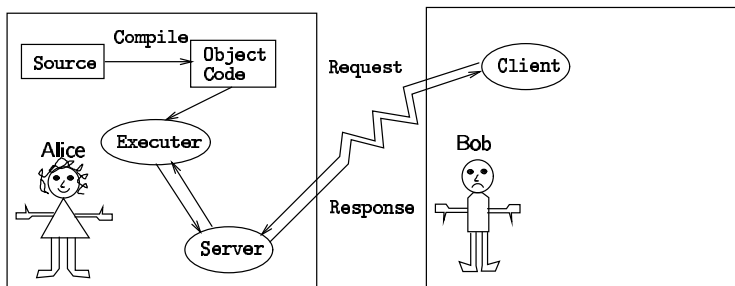
A rival developer (Bob) may steal proprietary algorithms or data structures.

Problem : it is becoming more common to distribute software in forms that are easy to decompile and reverse engineer.

Example : Java bytecode

Code obfuscation : Alice can protect her code by making reverse engineering so technically difficult that it becomes economically inviable.

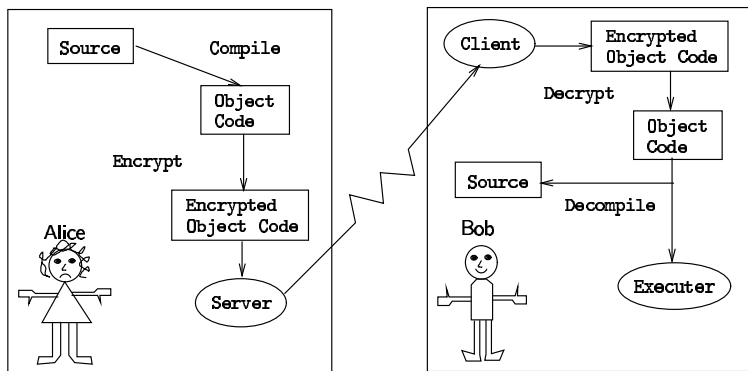
Other approach : Server-side protection



@C. Collberg

- ☺ Bob has not physical access to the code
- ☹ Bad performance

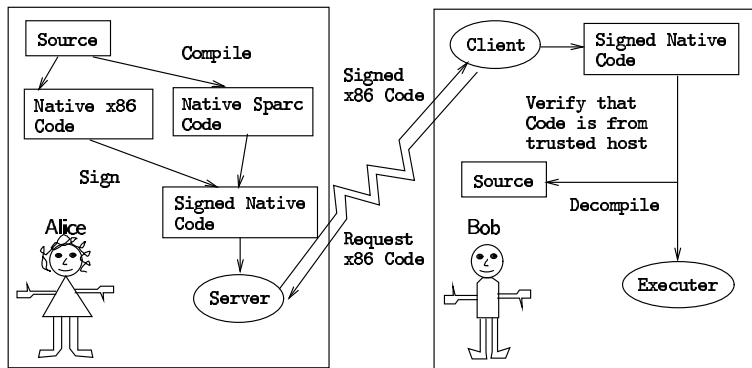
Other approach : protection by encryption



©C. Collberg

- ☺ Bob must decrypt the code to reverse engineer it
- ☹ decryption/execution must take place in hardware (\neq JVM)

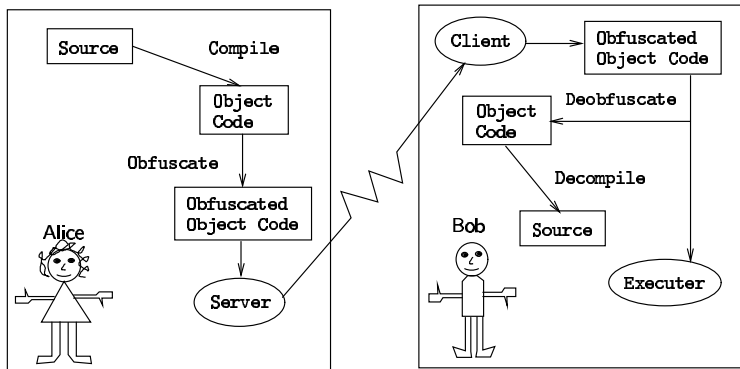
Other approach : protection by native code



©C. Collberg

- 😊 native code is more difficult to reverse engineer
- 😞 type safety is more difficult to ensure statically (unlike Java bytecodes)

Code obfuscation



©C. Collberg

- 😊 avoids previous shortcomings
- 😞 deobfuscation is always possible in theory (with enough time)

Obfuscation : definition

Let $P \xrightarrow{\mathcal{T}} P'$ be a transformation of a source program P into a target program P' . $P \xrightarrow{\mathcal{T}} P'$ is a *conservating transformation* if P and P' have the same *observable behavior* :

- 1 if P fails to terminate or terminates with an error, then P' may or may not terminate.
- 2 otherwise P' must terminate and produce the same output as P .

Quality

Collberg *et al.*² classify each transformations \mathcal{T} according to several criteria

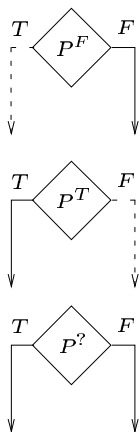
- potency** how much obscurity \mathcal{T} adds to the program
- resilience** how difficult \mathcal{T} is for a deobfuscator to undo,
- stealth** how well code introduced by \mathcal{T} fits in with the original code,
- cost** how much computational overhead \mathcal{T} adds to the obfuscated application.

2. C. Collberg, C. Thomborson, D. Low, *A Taxonomy of Obfuscating Transformations*

Trivial transformations

- ▶ remove the source code formatting informations
- ▶ remove comments (if any)
- ▶ scramble identifiers names

Opaque predicates



- The complexity of a program grows with the number of predicates it contains.
- A predicate P is opaque if its outcome is known at obfuscation time:
 - $P^F \Rightarrow P$ is always **False**.
 - $P^T \Rightarrow P$ is always **True**.
 - $P^? \Rightarrow P$ is sometimes **True**, sometimes **False**.
- Obfuscating control transformations insert opaque predicates that are difficult for a de-obfuscator to evaluate.

Simple opaque constructs

Trivial opaque constructs can be cracked by a deobfuscator by a static local analysis :

```
{ int v, a=5; b=6;
  v=11 = a + b;
  if (b > 5)T ...
  if (random(1,5) < 0)F ... }
```

©C. Collberg

Weak opaque constructs can be cracked by a deobfuscator by a static global analysis :

```
{ int v, a=5; b=6;
  if (...) ...
    ⋮ (b is unchanged)
  if (b < 7)T a++;
  v=36 = (a > 5)?v=b*b:v=b }
```

©C. Collberg

Numeric opaque constructs

Number theory textbooks may give inspiration...

$$\forall x, y \in \mathcal{I}, 7y^2 - 1 \neq x^2$$

$$\forall x \in \mathcal{I}, 2|(x + x^2)$$

$$\forall x \in \mathcal{I}, 3|(x^3 - x)$$

$$\forall n \in \mathcal{I}^+, x, y \in \mathcal{I}, (x - y)|(x^n - y^n)$$

$$\forall n \in \mathcal{I}^+, x, y \in \mathcal{I}, 2|n \vee (x + y)|(x^n + y^n)$$

$$\forall n \in \mathcal{I}^+, x, y \in \mathcal{I}, 2 \nmid n \vee (x + y)|(x^n - y^n)$$

$$\forall x \in \mathcal{I}^+, 9|(10^x + 3 \cdot 4^{(x+2)} + 5)$$

$$\forall x \in \mathcal{I}, 3|(7x - 5) \Rightarrow 9|(28x^2 - 13x - 5)$$

$$\forall x \in \mathcal{I}, 5|(2x - 1) \Rightarrow 25|(14x^2 - 19x - 19)$$

$$\forall x, y, z \in \mathcal{I}, (2 \nmid x \wedge 2 \nmid y) \Rightarrow x^2 + y^2 \neq z^2$$

$$\forall x \in \mathcal{I}^+, 14|(3 \cdot 7^{4x+2} + 5 \cdot 4^{2x-1} - 5)$$

Opaque predicates : other techniques

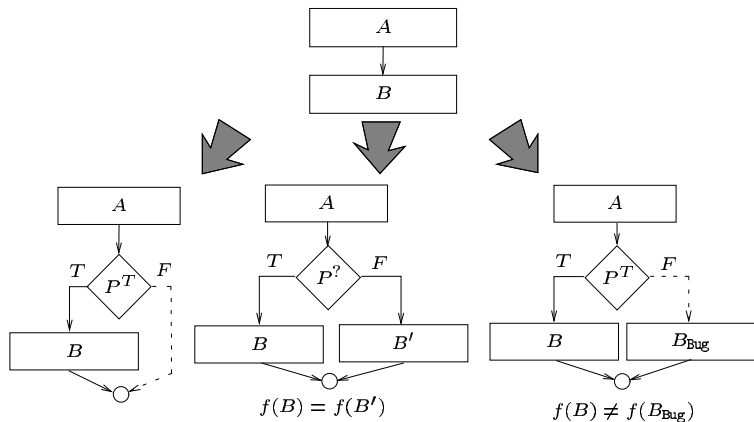
The other techniques target the weaknesses of the static analyses

- ▶ Alias : the program maintains complex alias invariants such that some predicates on pointer equality are constant.
- ▶ Concurrency : thread interleaving is difficult to predict statically

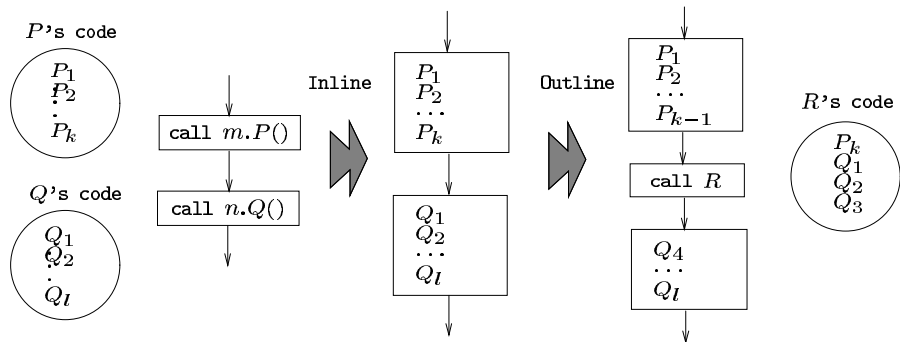
Using opaque predicates

Complexity (for the human understanding) increases with the number of predicates and nested levels of conditionals.

Application : insert Dead/Irrelevant code



Method inlining and outlining



©C. Collberg

Method interleaving

```

class C {
    A(T1 x){
        a1; a2;
    }
    B(T1 y; T2 z){
        b1; b2;
    }
}

{ C p=new C;
  p.A(x);
  p.B(y, z); }

class C' {
    M(T1 xy; T2 y; int V){
        if (V == p) {a1; a2;}
        else       {b1; b2;}
    }
}

{ C' p=new C';
  p.M(x, c, V=p);
  p.M(y, c, V=q); }

```

\xRightarrow{T}

Data transformations

Boolean variables and other variables of restricted range can be split into two or more variables.

(1) <code>bool A,B,C;</code>	(1') <code>short a1,a2,b1,b2,c1,c2;</code>
(2) <code>A = True;</code>	(2') <code>a1=0; a2=1;</code>
(3) <code>B = False;</code>	(3') <code>b1=0; b2=0;</code>
(4) <code>C = False;</code>	(4') <code>c1=1; c2=1;</code>
(5) <code>C = A & B;</code>	(5') <code>x=AND[2*a1+a2,2*b1+b2]; c1=x/2; c2=x%2;</code>
(6) <code>C = A & B;</code>	(6') <code>c1=(a1 ^ a2) & (b1 ^ b2); c2=0;</code>
(7) <code>C = A B;</code>	(7') <code>x=OR[2*a1+a2,2*b1+b2]; c1=x/2; c2=x%2;</code>
(8) <code>if (A) ...;</code>	(8') <code>x=2*a1+a2; if ((x==1) (x==2)) ...;</code>
(9) <code>if (B) ...;</code>	(9') <code>if (b1 ^ b2) ...;</code>
(10) <code>if (C) ...;</code>	(10') <code>if (VAL[c1,c2]) ...;</code>



String obfuscation

A constant string can be converted into a program that produces the string.

Exemple : G(1)="AAA" G(2)="BAAAA" G(3)=G(5)="CCB" G(4)="XCB"

```
String G (int n) {
    int i=0,k;
    String S;
    while (1) {
        L1:  if (n==1) {S[i++]="A";k=0;goto L6};
        L2:  if (n==2) {S[i++]="B";k=-2;goto L6};
        L3:  if (n==3) {S[i++]="C";goto L9};
        L4:  if (n==4) {S[i++]="X";goto L9};
        L5:  if (n==5) {S[i++]="C";goto L11};
            if (n>12) goto L1;
        L6:  if (k++<=2) {S[i++]="A";goto L6}
            else goto L8;
        L8:  return S;
        L9:  S[i++]="C"; goto L10;
        L10: S[i++]="B"; goto L8;
        L11: S[i++]="C"; goto L12;
        L12: goto L10;
    }
}
```

Array transformations

Arrays can be

- ▶ split into several sub-arrays
- ▶ merged : two or more arrays are merged into one
- ▶ folded : increase the number of dimension
- ▶ flattened : decrease the number of dimension

Array splitting

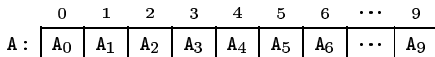
```
(1) int A[9];
```

```
(2) A[i] = ...;
```

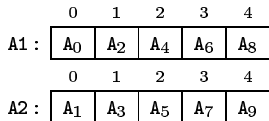
⇓ \mathcal{T}

```
(1') int A1[4], A2[4];
```

```
(2') if ((i%2)==0) A1[i/2]=...
      else A2[[i/2]]=...;
```



⇓ \mathcal{T}



Array splitting

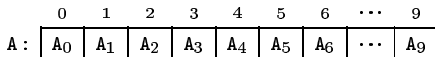
```
(1) int A[9];
```

```
(2) A[i] = ...;
```

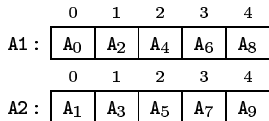
⇓ \mathcal{T}

```
(1') int A1[4], A2[4];
```

```
(2') if ((i%2)==0) A1[i/2]=...
      else A2[[i/2]]=...;
```



⇓ \mathcal{T}



Array merging

```
(3) int B[9],C[19];
```

```
(4) B[i] = ...;
```

```
(5) C[i] = ...;
```



```
(3') int BC[29];
```

```
(4') BC[3*i] = ...;
```

```
(5') BC[i/2*3+1+i%2] = ...;
```

	0	1	2	3	4	5	6	...	9
B :	B ₀	B ₁	B ₂	B ₃	B ₄	B ₅	B ₆	...	B ₉

	0	1	2	3	4	5	6	...	19
C :	C ₀	C ₁	C ₂	C ₃	C ₄	C ₅	C ₆	...	C ₁₉



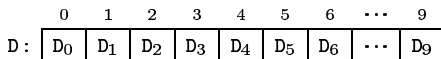
	0	1	2	3	4	5	6	...	29
BC :	B ₀	C ₀	C ₁	B ₁	C ₂	C ₃	B ₂	...	C ₁₉

Array folding

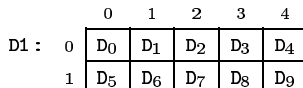
```
(6) int D[9];
(7) for(i=0;i<=8;i++)
    D[i]=2*D[i+1];
```

\Downarrow \mathcal{T}

```
(6') int D1[1,4];
(7') for(j=0;j<=1;j++)
    for(k=0;k<=4;k++)
        if (k==4)
            D1[j,k]=2*D1[j+1,0];
        else
            D1[j,k]=2*D1[j,k+1];
```



\Downarrow \mathcal{T}



Array flattening

```
(8) int E[2,2];
(9) for(i=0;i<=2;i++)
    for(j=0;j<=2;i++)
        swap(E[i,j], E[j,i]);
```

⇓ \mathcal{T}

```
(8') int E1[8];
(9') for(i=0;i<=8;i++)
    swap(E[i], E[3*(i/3)+i/3]);
```

	0	1	2
E : 0	E _{0,0}	E _{0,1}	E _{0,2}
1	E _{1,0}	E _{1,1}	E _{1,2}
2	E _{2,0}	E _{2,1}	E _{2,2}

⇓ \mathcal{T}

	0	1	2	3	4	...	8
E1 :	E _{0,0}	E _{0,1}	E _{0,2}	E _{1,0}	E _{1,1}	...	E _{2,2}

Conclusions

Theoretical limits

- ▶ Deobfuscation is always possible with enough time,
 - ▶ but some obfuscating transformations can be applied in polynomial time and require worst-case exponential time to deobfuscate (*i.e.* opaque predicates by aliasing).
- ▶ Obfuscation is also used by attackers (Bob) : viruses are duplicated and obfuscated to foil anti-virus software.
 - ▶ so deobfuscators are also necessary to ensure security !
- ▶ Obfuscation has received relatively little attention by the research community
 - ▶ but software companies use it
Example : gMail midlet